

EXPRESS MAIL MAILING LABEL NUMBER EV 318618392 US

PATENT
10830.0108.NPUS00

APPLICATION FOR UNITED STATES LETTERS PATENT

for

**DATABASE BLOCK NETWORKED ATTACHED STORAGE
PACKET JOINING**

by

John Ormond

BACKGROUND OF THE INVENTION

1. Field of the Invention

[0001] The present invention relates generally to data processing networks, and more particularly to database systems and network attached storage.

3. Description of Related Art

6 [0002] Database systems use block-level access to a database in data storage.

7 A database application issues block I/O requests to the database. Well-known database

8 applications include the IBM DB2, Oracle 8, and Sybase. The database applications may

9 support on-line transaction processing and batch processing.

10 [0003] Traditionally, the database application ran on a host processor such a
11 mainframe computer, and the database was stored in one or more disk drives directly
12 attached to the host processor. In the last couple of years, however, some database
13 systems have used network attached storage (NAS). For example, Celerra (Trademark)
14 brand network attached storage is sold by EMC Corporation, 176 South Street,
15 Hopkinton, MA 01748. There can be considerable savings in storage cost and data
16 management cost because multiple host processors can share the network attached
17 storage.

SUMMARY OF THE INVENTION

19 [0004] It has been discovered that many database applications using network
20 attached storage suffer a significant degradation in performance under high loading
21 conditions due to inefficient packing of block-level I/O requests into the network data
22 transport packets. For Transmission Control Protocol (TCP) transmission over an IP
23 network, for example, the IP network is often configured for a maximum transfer unit

1 (MTU) frame size of 9000 bytes, which is sufficient for transporting an 8 kilobyte data
2 block in each frame. When a host processor is concurrently executing multiple on-line
3 transaction processing (OLTP) application, many of the 9000 MTU frames will be less
4 than half full, and some of the frames will contain less than 500 bytes. There is a
5 considerable waste of host processing time and network bandwidth for transporting many
6 nearly empty frames.

7 [0005] In accordance with one aspect, the invention provides a method of
8 processing a series of data packets for transmission over a data network in a series of
9 frames in which at least some of the frames contain multiple data packets. Each data
10 packet in the series of data packets has a respective time in a time sequence. Each frame
11 is capable of transmitting a certain amount of data. The method includes successively
12 joining data packets from the time sequence into the frames and transmitting each data
13 packet in at least one of the frames no later than a certain time interval after the
14 respective time of said each data packet in the time sequence. The method also includes
15 transmitting each frame in a first set of the frames upon filling said each frame in the first
16 set of frames with data from one or more of the data packets so that said each frame in
17 the first set of frames cannot contain an additional data packet, and transmitting each
18 frame in a second set of the frames which are not filled with at least some of the data
19 packets so that said each frame in the second set of the frames cannot contain an
20 additional data packet in order to ensure that said each data packet is transmitted in at
21 least one of the frames no later than the certain time interval after the respective time of
22 said each data packet in the time sequence.

1 **[0006]** In accordance with another aspect, the invention provides a method of
2 operation in a host processor programmed for executing on-line transaction processing
3 applications and having a network block storage interface for accessing network attached
4 storage coupled to the host processor via a data network. The method includes the host
5 processor joining the I/O request data packets from different ones of the on-line
6 transaction processing applications in the same network transmission frames to more
7 completely fill the network transmission frames.

8 **[0007]** In accordance with yet another aspect, the invention provides a method
9 of solving a performance problem in a host processor programmed for executing on-line
10 transaction processing applications and having a network block storage interface for
11 accessing network attached storage coupled to the host processor via a data network. The
12 performance problem is caused by network transmission frames being only partially
13 filled with I/O request packets from the on-line transaction processing applications. The
14 performance problem is solved by re-programming the host processor to join the I/O
15 request data packets from different ones of the on-line transaction processing applications
16 in the same network transmission frames to more completely fill the network
17 transmission frames.

18 In accordance with a final aspect, the invention provides a host processor
19 programmed for executing on-line transaction processing applications and having a
20 network block storage interface for accessing network attached storage coupled to the
21 host processor via a data network. The host processor is programmed for joining the I/O
22 request data packets from different ones of the on-line transaction processing applications

1 into the same network transmission frames to more completely fill the network
2 transmission frames.

3 **BRIEF DESCRIPTION OF THE DRAWINGS**

4 **[0008]** Other objects and advantages of the invention will become apparent
5 upon reading the detailed description with reference to the drawings, in which:

6 **[0009]** FIG. 1 is a block diagram of a data processing system incorporating
7 the present invention;

8 **[00010]** FIG. 2 shows an inefficient method of packing I/O requests into
9 TCP/IP MTU frames;

10 **[00011]** FIG. 3 shows a more efficient method of packing I/O requests into
11 TCP/IP frames;

12 **[00012]** FIG. 4 shows various routines and data structures in an I/O request
13 bunching module introduced in FIG. 1;

14 **[00013]** FIG. 5 is a flowchart of an I/O request bunching main routine;

15 **[00014]** FIG. 6 is a flowchart of an I/O request bunching periodic timer
16 interrupt routine;

17 **[00015]** FIG. 7 is a flowchart of a procedure for turning on bunching of I/O
18 requests;

19 **[00016]** FIG. 8 is a flowchart of a procedure for turning off bunching of I/O
20 requests;

21 **[00017]** FIG. 9 is a flowchart of a procedure for configuration and adjustment
22 of a time interval "x" which, when exceeded, causes dumping of joined I/O requests to a
23 TCP/IP interface;

1 **[00018]** FIG. 10 shows various routines and data structures in an I/O request
2 bunching module that bunches read requests together and bunches write requests together
3 so that the order of the read requests and write requests may change;

4 **[00019]** FIG. 11 is a flowchart of read and write I/O request bunching in which
5 the order of the read requests and write requests may change;

6 **[00020]** FIGS. 12 and 13 comprise a flowchart of an I/O request bunching
7 main routine that that bunches read requests together and bunches write requests together
8 so that the order of the read requests and write requests may change;

9 **[00021]** FIG. 14 is a flowchart of an I/O request bunching timer interrupt
10 routine for use with the main routine of FIGS. 12 and 13; and

11 **[00022]** FIG. 15 shows a block diagram of I/O request bunching in a multi-
12 threaded system in which I/O request data packets in a range of I/O controller memory
13 are joined and packed in preallocated MTU frames.

14 **[00023]** While the invention is susceptible to various modifications and
15 alternative forms, specific embodiments thereof have been shown by way of example in
16 the drawings and will be described in detail. It should be understood, however, that it is
17 not intended to limit the form of the invention to the particular forms shown, but on the
18 contrary, the intention is to cover all modifications, equivalents, and alternatives falling
19 within the scope of the invention as defined by the appended claims.

20

21 **DESCRIPTION OF ILLUSTRATIVE EMBODIMENTS**

22 **[00024]** With reference to FIG. 1, there is shown a data processing system
23 incorporating the present invention. The data processing system includes a host

1 processor 21 coupled to multiple user terminals 22, 23, 24 for on-line transaction
2 processing. The host processor 21 is also coupled via an IP network 25 to network
3 attached storage 26.

4 **[00025]** The host processor includes on-line transaction processing applications
5 27 that send block-level I/O requests down to a network block storage TCP/IP interface
6 29. For example, the on-line transaction processing applications are separate instances of
7 a transaction processing program such as an accounting program for handling accounts
8 receivable. For each transaction, such as the processing of a single check, a customer's
9 account of money owed is debited by the amount of the check, and a vendor's account of
10 money received is credited by the amount of the check. The block-level I/O requests, for
11 example, are SCSI or SCSI-3 commands.

12 **[00026]** The Network Block Storage TCP/IP interface 29 receives data packets
13 from the on-line transaction processing applications 27, and each data packet includes
14 one or more block-level I/O requests. Upon receipt of a data packet, the Network Block
15 Storage TCP/IP interface places the data from the data packet in as many MTU frames as
16 required to hold all of the data of the data packet, and sends the MTU frames over the IP
17 network to the network attached storage 26.

18 **[00027]** The network attached storage 26 has a TCP/IP interface for removing
19 the block-level I/O requests from the MTU frames, and sending the block-level access
20 commands to a storage manager 32 that manages storage 35 containing a database 36.
21 The storage manager 32 maps logical block addresses referenced in the block-level I/O
22 requests to physical addresses in the storage 35.

1 **[00028]** It has been discovered that in an on-line transaction processing system
2 employing network storage as described above, there is a significant degradation in
3 performance under high loading conditions due to inefficient packing of block-level I/O
4 requests into the MTU frames. In particular, each of the on-line transaction processing
5 applications 27 may group a number of I/O requests together in a data packet before
6 sending the data packet down to the network block storage TCP/IP interface 29, but I/O
7 requests from one on-line transaction processing application are not grouped with another
8 on-line transaction processing application in a data packet. In contrast to an off-line or
9 batch transaction processing application, a majority of the data packets from an on-line
10 transaction processing application may have a relatively small size compared to the MTU
11 frame size. For example, the data packets are often only 500 bytes, and the MTU frame
12 size is typically configured as either 1,500 bytes or 9,000 bytes.

13 **[00029]** FIG. 2 for example, shows a series of I/O request packets 41, 42, 43,
14 and 44 for a system in which a host processor is concurrently executing multiple on-line
15 transaction processing (OLTP) applications. The I/O request packets 41, 42, 43, and 44
16 are transmitted in respective TCP/IP MTU frames 51, 52, 53, 54. In this system, each I/O
17 request packet may include one or more requests, but each I/O request packet originates
18 from one of the applications. Moreover, each of the I/O request packets is placed in a
19 respective one of the MTU frames. For I/O request packets that are small in comparison
20 to the MTU frame size, there is a considerable waste of host processing time and network
21 bandwidth for transporting many nearly empty frames.

22 **[00030]** The problem introduced in FIG. 2 can be solved by joining I/O
23 requests as shown in FIG. 3. A number of the I/O request packets 41, 42, and 43 are

1 joined together and placed in the same MTU frame 55 until no more of the I/O request
2 packets can be placed in the MTU frame. As described further below, such joining of the
3 data packets into an MTU frame and initiation of transmission of the MTU frame can be
4 performed by a main routine 61 in FIG. 4. Therefore, the frame 55 becomes filled with
5 data packets 41, 42, 43 so that the frame contain the additional data packet 44. Some
6 delay is introduced in the transmission of most of the I/O request packets, but this delay
7 can be limited to a certain value “x”. For example, the I/O request 44 is transmitted alone
8 in the MTU frame 56 to ensure that the delay is limited to the certain value “x”. As
9 described further below, the transmission of such frames to satisfy the delay constraint
10 can be initiated by a timer interrupt routine 62 in FIG. 4.)

11 [00031] In the host processor 21 of FIG. 1, each of the on-line transaction
12 processing applications 27 sends block level I/O requests to the I/O request bunching
13 module 28. When the I/O request bunching module 28 receives an I/O request that is
14 smaller than the MTU frame size, the I/O request bunching module will attempt to join
15 the I/O request with another block-level I/O request in order to more completely fill an
16 MTU frame. However, the I/O request bunching module 28 will not withhold an I/O
17 request from the network block storage TCP/IP interface 29 for more than a certain time
18 interval “x”. When the network block storage TCP/IP interface 29 receives an I/O
19 request or joined I/O requests, it packs the I/O request or the joined I/O requests in as
20 many MTU frames as needed and transmits the MTU frames over the IP network 25 to
21 the network attached storage 26. The network block storage TCP/IP interface is
22 configured to use the 9000 byte MTU frame size. In the future, if the TCP/IP interface

1 could be configured to use an MTU frame size larger than 9000 bytes, then it may be
2 desirable to use the larger MTU frame size.

3 [00032] In the network attached storage 26, the TCP/IP interface strips the I/O
4 requests from the MTU frames, and sends the I/O requests to the storage manager 32.
5 The storage manager 32 interprets the I/O requests for read and write access to the
6 database 36 in storage 35, and formulates a reply to each I/O request. The replies are
7 received in an I/O reply bunching module 33. The I/O reply bunching module 33
8 functions in a similar fashion as the I/O request bunching module 28. In particular, when
9 the I/O reply bunching module 33 receives an I/O reply that is smaller than the MTU
10 frame size, the I/O reply bunching module will attempt to join the I/O reply with another
11 I/O reply in order to more completely fill an MTU frame. However, the I/O reply
12 bunching module will not withhold an I/O reply from the network block storage TCP/IP
13 interface 29 for more than a certain time interval "x", which can be the same time interval
14 used by the I/O request bunching module.

15 [00033] For request or reply bunching, the time interval "x" can be initially
16 selected based on application type. For example, the time interval "x" can be set as a
17 fraction of the nominal I/O response time for the on-line transaction processing
18 application. The nominal I/O response time is an average I/O response time of the
19 application when request and reply bunching is not used and the IP network is lightly
20 loaded. In particular, for an application having a nominal I/O response time of 15
21 milliseconds, the value of "x" can be 5 milliseconds.

22 [00034] It is also desirable to disable I/O request bunching if large bursty
23 transactions occur so that most of the MTU frames would become nearly full without

1 joining I/O requests or replies. Such large bursty transactions occur in Datawarehouse
2 databases where bulk database data is being moved to and from the database. A
3 bunching flag can be set and cleared to enable and disable the bunching of I/O requests
4 and replies. This bunching flag could be cleared when the bulk database transfers occur
5 and otherwise set. The bunching flag could also be cleared during certain times of the
6 day when bulk database transfers are likely to occur.

7 [00035] The request bunching module 28 in FIG. 1 can be constructed in
8 various ways. One way is shown in FIG. 4. The request module 28 includes a main
9 routine 61, a timer interrupt routine 62, a request joining buffer 63, a joined size variable
10 64, an oldest request time variable 65, and a bunching flag 66. The I/O reply bunching
11 module (33 in FIG. 1) can be constructed in the same fashion.

12 [00036] The request bunching module 28 in FIG. 1 could be programmed to
13 function in various ways. One way is for the request bunching module to be an add-in
14 program that intercepts I/O request data packets sent from the on-line transaction
15 processing applications to the network block storage TCP/IP interface. In the absence of
16 the I/O request bunching module 28, the network block storage TCP/IP interface would
17 receive each I/O request data packet, pack the respective I/O request data packet into one
18 or more MTU frames, and then transmit the MTU frames over the IP network 25.

19 [00037] During bunching, the I/O request bunching module 28 receives a series
20 of consecutive input I/O request data packets from the on-line transaction processing
21 applications 27, joins the respective I/O request data packets from the consecutive I/O
22 requests to form a joined I/O request data packet, and then transmits the joined I/O
23 request data packet to the network block storage TCP/IP interface 29. The process of

1 joining consecutive I/O requests to form the joined I/O request data packet will terminate
2 with the transmission of the joined I/O data packet to the network block storage TCP/IP
3 interface once the joined data packet reaches the MTU frame data block size (e.g., 8
4 kilobytes for a 9000 MTU frame) or when needed to ensure that transmission of the
5 joined I/O data packet is not delayed by more than the time interval "x".

6 [00038] When the I/O request bunching module is constructed as an add-in
7 program, it is desirable for the request bunching module to use very few host processor
8 execution cycles for processing each I/O request data packet intercepted from the on-line
9 transaction processing applications, and to use very few host processor execution cycles
10 for transmitting each joined I/O request to the network block storage TCP/IP interface 29.
11 For example, a timer interrupt routine 62 separate from the main routine 61 is used to
12 check whether the time interval "x" has expired instead of programming the main routine
13 61 to check for expiration of the time interval "x" each time that the main routine
14 processes an I/O request intercepted from the on-line transaction processing applications
15 27. This significantly reduces the number of host processor execution cycles used when
16 processing a large number of small I/O request data packets over a short interval of time.

17 [00039] FIG. 5 shows a flowchart for the main routine of the I/O request
18 bunching module. This routine is started when one of the on-line transaction processing
19 applications sends a new I/O request data packet to the I/O request bunching module. In
20 a first step 71, if the bunching flag is not set, then execution branches to step 72 to pass
21 the new I/O request data packet to the network block storage TCP/IP interface, and the
22 main routine is finished.

1 **[00040]** In step 71, if the bunching flag is set, then execution continues from
2 step 71 to step 73. In step 73, the main routine gets a lock on the request joining buffer.
3 This lock is used to avoid conflict since access to the request joining buffer is shared at
4 least with the timer interrupt routine. In step 74, if the buffer is empty, then execution
5 continues to step 75. In step 75, the oldest request time variable is set to the present time.
6 Execution continues from step 75 to step 76. Execution also continues to step 76 from
7 step 74 when the request joining buffer is not empty.

8 **[00041]** In step 76, the main routine inserts the new I/O request data packet
9 onto the tail of the request joining buffer. In step 77, the main routine increments the size
10 of the joined request (i.e., the joined size variable) by the size of the new I/O request data
11 packet. In step 78, if the joined size is greater or equal to 8 kilobytes, then execution
12 continues to step 78. In step 79, the joined request is passed to the NBS TCP/IP
13 interface. In step 80, the request joining buffer is cleared. In step 81, the time of the
14 oldest request is cleared. (For the comparison in step 92 as described below, the time of
15 the oldest request should be cleared by setting it to a very high value, always representing
16 a time in the future, so that the timer interrupt routine effectively does nothing until the
17 time of the oldest request becomes set to the present time in step 75.) In step 82, the lock
18 on the request joining buffer is released, and the main routine is finished.

19 **[00042]** In step 78, if the joined size is not greater than or equal to 8 kilobytes,
20 then execution continues to 82 to release the lock, and the main routine is finished.

21 **[00043]** FIG. 6 shows the timer interrupt routine 62. This routine is executed
22 periodically, for example, once every millisecond. In a first step 91, if the request joining
23 buffer is empty, then the timer interrupt routine is finished. Otherwise, execution

1 continues from step 91 to step 92. In step 92, if the difference between the present time
2 and the value of the oldest request time variable is not greater than "x1", then execution
3 returns. Otherwise, the time interval "x1" has been exceeded, and execution continues
4 from step 92 to step 93. "x1" is the time interval "x" referred to above minus the period
5 of the periodic interrupt of FIG. 6.

6 [00044] In step 93, the periodic timer interrupt routine gets a lock on the
7 request joining buffer. In step 94, the joined request is passed to the network block
8 storage TCP/IP interface. In step 95, the request joining buffer is cleared. In step 96 the
9 oldest request time is cleared, for example, by setting it to a very high value representing
10 a time that is always in the future. In step 97, the lock on the request joining buffer is
11 released, and the timer interrupt routine is finished.

12 [00045] Bunching can be turned on and off dynamically. For example, FIG. 7
13 shows a procedure for turning on bunching. In a first step 101, the periodic timer
14 interrupt routine is enabled. In step 102, the bunching flag is set, and the procedure is
15 finished.

16 [00046] FIG. 8 shows a procedure for tuning off bunching. Bunching is turned
17 off in such a way that the request joining buffer is flushed and the timer interrupt routine
18 is disabled. In a first step 111, a lock is obtained on the request joining buffer. In step
19 112, if the buffer is not empty, execution continues to step 113. In step 113, the joined
20 I/O request data packet is passed to the network block server TCP/IP interface. In step
21 114, the request joining buffer is cleared. Execution continues from step 114 to step 115.
22 Execution also continues to step 115 from step 112 when the request joining buffer is
23 empty.

1 **[00047]** In step 115, the oldest request time is cleared, for example, by setting
2 it to a very high value representing a time that is always in the future. In step 116, the
3 request bunching flag is cleared. In step 117, the lock on the request joining buffer is
4 released. In step 118, the timer interrupt routine is disabled, and the procedure is
5 finished.

6 **[00048]** The time interval "x" can also be dynamically adjusted based on
7 loading characteristics of the IP network. For example, the data processing system of
8 FIG. 1 has a load monitor 37 that measures loading of the IP network with respect to the
9 handling of the I/O request and replies, and uses the measured loading to adjust the time
10 interval "x". For example, the measured loading ranges from zero for no loading, to one
11 for saturation of the IP network for the transmission of the I/O requests and replies. The
12 time interval "x" can be adjusted based on a formula of the loading such as:

13
$$x = x_{\min} + (x_{\max} - x_{\min})(\text{loading})$$

14 **[00049]** In this example, the time interval "x" ranges from a minimum of x_{\min}
15 for a loading of zero to a maximum of x_{\max} for a loading of one. The value of x_{\max} can
16 be set to a fraction of the nominal I/O response time, and the value of x_{\min} can be set to a
17 convenient minimum time for the checking for joined requests over the interval "x". In
18 particular, for an application having a nominal I/O response time of 15 milliseconds, the
19 value of " x_{\max} " can be 7 milliseconds, and the value of " x_{\min} " can be 1 or 2 milliseconds.

20 **[00050]** The value of "x" could also be adjusted or based on the average size of
21 the data blocks in the database 35 or the average size of the I/O request packets as
22 received by the I/O request bunching. For example, for larger I/O request packet size, a
23 smaller size of "x" could be used. If the I/O request packet size is substantially different

1 from the I/O reply packet size, then it may be desirable for the size of "x" for the I/O
2 request bunching 28 to be different from the size of "x" for the I/O reply bunching.

3 [00051] The value of "x" can be adjusted dynamically based on data activity or
4 caching algorithms in the host processor 21 and in the network attached storage 26 in
5 addition to the monitored loading on the IP network. Estimated loading on the IP
6 network could also take into account activity of any other hosts that may share the IP
7 network 25. The average size of a data block in the database 35 could be dynamically
8 calculated in the applications 27 and passed down to the I/O request bunching module for
9 adjustment of "x". The average size of the I/O request packets as received by the I/O
10 request bunching could be dynamically calculated in the I/O request bunching module
11 itself.

12 [00052] The I/O request bunching module could also estimate its loading on
13 the IP network by accumulating the average number of blocks of joined requests per unit
14 time and the average number of bytes in the joined requests per unit time, computing an
15 estimate of the loading as a function of the average number of blocks of joined requests
16 per unit time and the average number of bytes in the joined requests per unit time, and
17 adjusting "x" based on the computed estimate of the loading.

18 [00053] In short, the value of "x" can be continually reset to achieve the best
19 performance in database access based on the current processing environment and I/O
20 activity of the host processor 21 and the network attached storage 26. The end result is a
21 more intelligent NAS transfer mechanism in which an optimum amount of small database
22 blocks in the I/O request and replies are accumulated and packed into the MTU frames.

1 **[00054]** FIG. 9 summarizes the configuration and dynamic adjustment of the
2 time interval "x". In a first step 121, the nominal I/O response time for the on-line
3 transaction processing application is determined. For example, the nominal I/O response
4 time is the average response time when request and reply bunching are not used and the
5 IP network is lightly loaded. In step 122, the time interval "x" is set to a fraction such as
6 one-third of the nominal I/O response time. In step 123, I/O request and reply bunching
7 is begun on the system while continually measuring the loading on the network and the
8 average size of the data blocks in the data base and the average size of the request packets
9 received by the I/O request, and dynamically adjusting the time interval "x" based on the
10 measurements. In step 124, the I/O request and reply bunching is disabled during large
11 bursty transactions such as bulk database data transfer.

12 **[00055]** In many systems, there is no need to preserve the ordering of read I/O
13 requests with respect to write I/O requests as the I/O requests are transmitted from the on-
14 line transaction processing applications to the network attached storage. In this case, it
15 may be possible to improve performance by separately bunching the read I/O requests
16 together, separately bunching the write I/O requests together, and dumping the bunched
17 read I/O requests before dumping the bunched write I/O requests when the time interval
18 "x" is exceeded. This improves performance because the I/O response time for reads is
19 generally faster than the I/O response time for writes. In addition, the bunched read
20 requests and the bunched write requests may tend to access separate localized regions of
21 memory so that there is a performance gain due to more efficient data caching and less
22 frequent read/write head arm swings in disk drives that comprise the storage (35 in FIG.
23 1) containing the database (36 in FIG. 1).

1 **[00056]** FIG. 10 shows an I/O request bunching module 130 for separately
2 bunching the read I/O request together, and separately bunching the write I/O requests
3 together. The I/O request bunching module 130 includes a main routine 131, a timer
4 interrupt routine 132, a read request joining buffer 133, a write request joining buffer
5 134, a read joined size variable 135, a write joined size variable 136, an oldest request
6 time variable 137, and a bunching flag 138.

7 **[00057]** FIG. 11 shows the overall operation of the request bunching module
8 130 of FIG. 10. In a first step 139 of FIG. 11, the read requests are bunched together in
9 the read request joining buffer, and the write requests are bunched in the write joining
10 buffer, so that the order of the read requests and write requests may be changed. In step
11 140, when the time interval "x1" is exceeded, or in response to a transaction commit
12 request from one of the on-line transaction processing applications, the read request
13 joining buffer is dumped to the TCP/IP interface, and then the write request joining
14 buffer is dumped to the TCP/IP interface. This dumping of the read requests before the
15 write requests gives priority to reads over writes, by moving some of the read I/O request
16 data packets in front of some of the write I/O request data packets in some of the frames.

17 **[00058]** FIGS. 12 and 13 show the main routine (131 in FIG. 10) for read and
18 write I/O request bunching. This main routine is begun upon receipt of a read or write
19 I/O request from one of the on-line transaction processing applications. In a first step 141
20 of FIG. 12, if the bunching flag is set, then execution branches to step 142 to pass the
21 new I/O request to the network block storage TCP/IP interface. Otherwise, if the
22 bunching flag is set, execution continues from step 141 to step 143. In step 143, if the
23 new request is a read request, then execution continues to step 144. In step 144, a lock is

1 obtained on the read request joining buffer. In step 145, if the read request joining buffer
2 is empty, then execution continues to step 146. In step 146, the oldest read request time
3 is set to the present time. Execution continues from step 146 to step 147 in FIG. 13.
4 Execution also branches to step 147 in FIG. 13 from step 146 if the read request joining
5 buffer is not empty.

6 [00059] In step 147 of FIG. 13, the new read request data packet is inserted
7 onto the tail of the read request joining buffer. In step 148, the size of the joined read
8 request data packet is incremented by the size of the new read request data packet. In
9 step 149, if the read joined size is greater than or equal to 8 kilobytes, then execution
10 continues to step 150. In step 150, the read request joining buffer is dumped to the
11 network block storage TCP/IP interface, for example, the contents of the read request
12 joining buffer are passed to the network block storage TCP/IP interface, the read request
13 joining buffer is cleared, and the time of the oldest read request is cleared. Execution
14 continues from step 150 to step 151. Execution also branches from step 149 to step 151
15 when the read joined size is not greater than or equal to 8 kilobytes. In step 151, the lock
16 on the read request joining buffer is released, and the main routine is finished.

17 [00060] In step 146 of FIG. 11, if the new request is a not a read request, then it
18 is a write request, and execution branches to step 152. In step 152, a lock is obtained on
19 the write request joining buffer. In step 153, if the write request joining buffer is empty,
20 then execution continues to step 154. In step 154, the oldest write request time is set to
21 the present time. Execution continues from step 154 to step 155 in FIG. 13. Execution
22 also branches to step 155 in FIG. 13 from step 153 if the read request joining buffer is not
23 empty.

1 **[00061]** In step 155 of FIG. 13, the new write request data packet is inserted
2 onto the tail of the write request joining buffer. In step 156, the size of the joined write
3 request data packet is incremented by the size of the new write request data packet. In
4 step 157, if the write joined size is greater than or equal to 8 kilobytes, then execution
5 continues to step 158. In step 158, the write request joining buffer is dumped to the
6 network block storage TCP/IP interface, for example, the contents of the write request
7 joining buffer are passed to the network block storage TCP/IP interface, the write request
8 joining buffer is cleared, and the time of the oldest write request is cleared. Execution
9 continues from step 158 to step 159. Execution also branches from step 157 to step 159
10 when the read joined size is not greater than or equal to 8 kilobytes. In step 159, the lock
11 on the write request joining buffer is released, and the main routine is finished.

12 **[00062]** FIG. 14 shows the periodic timer interrupt routine 132 for read and
13 write I/O request bunching. In a first step 161, if the read and write request joining
14 buffers are empty, then execution returns. Otherwise, execution continues from step 161
15 to step 162. In step 162, if the present time minus the oldest request time (i.e., the oldest
16 of the oldest read request time and the oldest write request time) is not greater than “x1”,
17 then execution returns. Otherwise, execution continues from step 162 to step 163. In
18 step 163, a lock is obtained on the read and write request joining buffers. Then in step
19 164, if the read request joining buffer is not empty, execution continues to step 165. In
20 step 165, the joined read request in the read request joining buffer is passed to the
21 network block services TCP/IP interface. Execution continues from step 165 to step 166.
22 Also, execution branches from step 164 to step 166 if the read request joining buffer is
23 empty.

1 **[00063]** In step 166, if the write request joining buffer is empty, then execution
2 continues from step 166 to step 167. In step 167, the joined write request in the write
3 request joining buffer is passed to the network block services TCP/IP interface.
4 Execution continues from step 167 to step 168. Also, execution branches from step 166
5 to step 168 if the write request joining buffer is empty.

6 **[00064]** In step 168, the read and write request joining buffers are cleared.
7 Then in step 169, the oldest request time is cleared, for example, by setting it to a large
8 value representing a time that is always in the future. Finally, in step 170, the lock is
9 released.

10 **[00065]** FIG. 15 shows a block diagram of I/O request bunching in a multi-
11 threaded system. Instead of performing I/O request bunching in an add-in function, the
12 request bunching occurs in the host processor's network block storage TCP/IP interface
13 during the transfer of I/O request data packets from the on-line transaction processing
14 applications to preallocated MTU frames.

15 **[00066]** As shown in FIG. 15, I/O controller memory 181 contains a pool 182
16 of a certain number "y" of preallocated MTU frames, and the pool 182 is updated every
17 "n" milliseconds. For example, the number of preallocated frames "y" is 400, and "n" is
18 two milliseconds for a system where the time interval "x" is 5 milliseconds. The I/O
19 controller memory 182 also has "r" ranges of addresses, which have received I/O request
20 data packets received within the "n" millisecond update interval. Each address range
21 contains the number of I/O requests that generally can be joined together and packed into
22 MTU frames by a single thread during the "n" millisecond update interval. In effect,

1 these address ranges function as an interface queue between the on-line transaction
2 processing applications and the network block storage TCP/IP interface.

3 [00067] Application threads 184 of the on-line transaction processing
4 applications load the I/O request data packets into the “r” address ranges of the I/O
5 controller memory 182. TCP/IP interface threads 185 preallocate MTU frames. In
6 addition, for each of the “r” address ranges, a respective one of the TCP/IP interface
7 threads 185 picks up I/O request data packets from the address range, joins these I/O
8 request data packets, and packs the joined I/O request data packets into preallocated
9 MTU frames. For example, a single thread generally processes a certain number “NPT”
10 of I/O requests during one “n” millisecond interval, and if there are “NR” new requests in
11 the “n” millisecond interval, then there are about $r = NR/NPT$ address ranges, and a
12 respective thread is initiated to process the I/O requests in each of the “r” address ranges.
13 TCP/IP threads 186 transmit the packed MTU frames over the IP network 25.

14 [00068] For some on-line transaction processing applications, the I/O requests
15 data packets are roughly the same size and are limited to a certain size “Z”. In this case,
16 each of the “r” address ranges can be the same size of “Z” times “NPT”.

17 [00069] The following table shows a specific example of a series of I/O
18 requests, the threads that concurrently join the I/O requests, and the MTU frames into
19 which the joined requests are packed:

Memory	Address	IO Request #	Size	IO Request Join Review	Applied Concurrency to large I/O Base	MTU Frame #
--------	---------	--------------	------	------------------------	---------------------------------------	-------------

000000	1	0.5k	Thread1	1
000001	2	1k		
000002	3	0.2k		
000003	4	0.3k		
000004	5	2k		
000005	6	4k		
000006	7	0.4k		2
000007	8	0.66k		
000008	9	0.45k		
000009	10	2.5k		
000010	11	.34k		
000011	12	1.6k		
000012	13	3.4k		
000013	14	10k		
000014	15	7k		
000015	16	.543k		
	etc			
			end of	
			range or	
	100		"r"	
000100	101	5k	Thread 2	3
000101	102	2k		
000102	103	0.1k		
000103	104	0.55k		4
000104	105	1k		
000105	106	.44k		
000106	107	0.67k		
000107	108	0.99k		
000108	109	3.5k		5

000109	110	6.7k	
000110	111	0.04k	
000111	112	1.2k	6
000112	113	5.2k	
000113	114	0.52k	7
	etc..	...	
		end of range or	
	200	"r"	

1

2

3 **[00070]** In the example shown by the above table, the I/O request data packets
 4 are organized as variable-length records in the controller memory, and the records are
 5 mapped to sequential record numbers that serve as respective memory addresses for the
 6 I/O request data packets. A first range of memory addresses from 000000 to 000099
 7 stores the I/O request data packets for a first set of one-hundred I/O requests. A first
 8 thread has joined and has packed the first six I/O request data packets into a first MTU
 9 frame. The first thread is presently joining the seventh I/O request data packet with the
 10 eighth I/O request data packet into a second MTU frame. Concurrently, a second thread
 11 has packed I/O request data packets 101, 102, and 103 into a third MTU frame; I/O
 12 request data packets 104 to 108 into a fourth MTU frame; I/O request data packets 109 to
 13 111 into a fifth MTU frame; and I/O request data packets 112 to 113 into a sixth MTU
 14 frame. The second thread is presently joining I/O request data packet 114 with following
 15 data packets into a seventh MTU frame.

1 [00071] The following table further shows fifteen threads concurrently
2 processing fifteen-hundred I/O request data packets received by the network block
3 storage TCP/IP interface within a one-millisecond interval:

4

"r" Memory Range	IO Request #	Thread
0001-000100	1-100	Thread1
000101-000200	101-200	Thread2
000201-000300	201-300	Thread3
000301-000400	301-400	Thread4
000401-000500	401-500	Thread5
000501-000600	501-600	Thread6
000601-000700	601-700	Thread7
000701-000800	701-800	Thread8
000801-000900	801-900	Thread9
000901-0001000	901-1000	Thread10
0001001-0001100	1001-1100	Thread11
0001101-0001200	1101-1200	Thread12
0001201-0001300	1201-1300	Thread13
0001301-0001400	1301-1400	Thread14
0001401-0001500	1401-1500	Thread15

5

6 [00072] The arrangement in FIG. 15 permits the TCP/IP interface to be
7 concurrently working on up to "y" MTU frames. The MTU frames can always be
8 preallocated and waiting to receive I/O request data packets by the time that they are
9 packed with the I/O request data packets. Therefore, once an MTU frame is filled with
10 the I/O request data packets or when the time interval "x" is exceeded, the MTU frame is
11 ready to be transmitted over the IP network 25. This will maintain a consistent amount of

1 traffic on the network pipe and will use a predictable amount of processing time on the
2 host processor. The overall performance is improved by more completely filling the MTU
3 frames with the I/O request data packets.

4 [00073] In view of the above, there has been described a performance problem
5 caused by network transmission frames being only partially filled with I/O request
6 packets from the on-line transaction processing applications. This performance problem
7 is solved by re-programming the host processor to join the I/O request data packets from
8 different ones of the on-line transaction processing applications in the same network
9 transmission frames to more completely fill the network transmission frames. Preferably
10 this is done by successively joining the I/O request data packets into the frames and
11 transmitting each data packet in a frame after a delay of no more than a certain time
12 interval. At least some of the frames are transmitted once these frames are filled with
13 some of the data packets so that each of these frames cannot contain an additional data
14 packet. Preferably the certain time interval is adjusted based on network loading so that
15 the certain time interval is increased for increased loading.

16